# Providing Access to Census-based Interaction Data in the UK: That's WICID!

John Stillwell

School of Geography, University of Leeds

Leeds, LS2 9JT, United Kingdom

## ABSTRACT

The Census Interaction Data Service (CIDS) is funded by the Economic and Social Research Council in the UK to provide access for social science researchers and students to the detailed migration and journey-to-work statistics that are collected by the national statistical agencies. These interaction data sets are known collectively as the Special Migration Statistics (SMS) and the Special Workplace Statistics (SWS). This paper outlines how problems of user access to these data have been tackled through the development of a web-based system known as WICID (Web-based Interface to Census Interaction Data). The paper illustrates various interface features including some of the query building facilities that enable users to extract counts of flows of particular groups of individuals between selected origin and destination areas. New tools are outlined for assisting area selection using digital maps of census geographies, for planning output and for adding value to the data through analysis. Mapping of flows of migrants between London boroughs and the rest of the UK demonstrates the value of the data. The paper begins with a summary of the data sets that are contained within the system and an outline of the system architecture.

Keywords: Census, interaction, geography, data, migration, commuting, interface

## 1. UK CENSUS INTERACTION DATA

Since 2002, the Census Interaction Data Service (CIDS) has been providing social science researchers and students in the UK with access to the 'Origin-Destination Statistics' that are collected by the UK statistical agencies (Office of National Statistics for England and Wales, General Register Office for Scotland, Northern Ireland Statistics and Research Agency) through decadal censuses. Origin-Destination Statistics, so-called because they involve two geographical areas, comprise the counts of migrants in the 12 months prior to each population census (in 2000-01, 1990-91 and 1980-81) and counts of those commuting to work at the time of the census (in 2001, 1991 and 1981). These data sets are known respectively as the Special Migration Statistics (SMS) and the Special Workplace Statistics (SWS), each of which is comprised of a set of tables containing a number of variables or data counts. The numbers of tables and counts available depend on geographical scale and there are essentially three levels of population census geography. Table 1 is a summary of the tables and counts that are available at each level, where level 1 refers to local authority districts, level 2 refers to census wards and level 3 refers to census output areas. In 2001, there were 426 districts, 10,608 wards and 223,060 output areas in the UK. The origin-destination matrices of flows thus contain a very large number of cells, particularly at the smallest spatial scale.

**Table 1: Tables and counts in the 2001 and 1991 interaction data sets for the UK**

| Data sets | Level 1: Districts |
|---|---|
| 2001 SMS | 10 tables, 996 counts |
| 1991 SMS | Set 2: 11 tables, 94 counts |
| 2001 SWS | 7 tables, 936 counts |
| 2001 STS | 7 tables, 1,176 counts |
| 1991 SWS* | - |
| | **Level 2: Wards** |
| 2001 SMS | 5 tables, 96 counts |
| 1991 SMS | Set 1: 2 tables, 12 counts |
| 2001 SWS | 6 tables, 354 counts |
| 2001 STS | 6 tables, 478 counts |
| 1991 SWS* | Set C: 9 tables, 274 counts |
| | **Level 3: Output areas** |
| 2001 SMS | 1 table, 12 counts |
| 1991 SMS | - |
| 2001 SWS | 1 table, 36 counts |
| 2001 STS | 1 table, 50 counts |
| 1991 SWS* | - |

*\* 10% sample*

Table 1 compares the volumes of data available from the last two censuses, indicating how more detailed data have become available from the latest census and demonstrating that more information is available for larger geographical areas (districts) than for small areas (output areas). The 2001 census was the first to have coordinated outputs for England and Wales, Scotland and Northern Ireland. However, GRO(S) decided to collect and produce data on journeys to study as well as journeys to work and hence Special Travel Statistics (STS) replace SWS for 2001 in Scotland.

These migration and commuting data sets are a hugely important resource for research and planning since there is no population registration system in the UK and there are no alternative sources that provide data of similar reliability or of equivalent spatial coverage. They are large and complex data sets with which the user needs to gain some familiarity before using them with confidence. The interaction data from the 1991 census have been reviewed by Flowerdew and Green (1993)[3] and shortcomings have been documented in detail by Rees *et al.* (2002)[4] and Cole *et al.* (2002)[2]. Some of the problems have arisen because of the requirement for the census agencies to preserve confidentiality and to lower the risk of disclosure. Various methods of doing this have been used in different censuses. In 1991, the Office of Population Censuses and Surveys (OPCS) suppressed flows of under 10 persons in many of the interaction data tables thereby creating significant undercounts of the flows taking place. To counter this problem, Rees and Duke-Williams (1997)[5] re-estimated the missing data to provide tables with sets of data derived from the primary counts. In 2001, the approach of ONS to reduce the risk of disclosure was a small cell adjustment method (SCAM), in

which values of 1 and 2 in the tables were seemingly adjusted to values of either 0 or 3.

Another problem with the census interaction data sets is that geographical boundaries change from year to year, thus obscuring the easy comparison of data over time. These changes tend to be more significant for smaller geographical areas such as wards than for large units such as districts. The smallest units, output areas, only appeared in the 2001 results since, for the first time, the geography used for data collection (enumeration districts) was not used for census outputs. The results of boundary change have prompted the CIDS team to produce estimates of flows from the 1991 and 1981 censuses for 2001 geographical areas. These derived data sets have been generated using a modelling methodology outlined in Boyle and Feng (2002)[1] and are available for users to extract from WICID along with all the primary data sets and the 1991 re-estimated data sets.

## 2. The WICID system

Because of the size and relative complexity of these data sets, they remain underused. Consequently, one of the challenges is to develop a user interface to them that will improve their accessibility and encourage their extraction and use. WICID is the web-based information system that has been developed to facilitate access to these primary interaction data sets and to additional sets of data derived from them. WICID has been designed to offer a user interface via a standard web browser; it is constructed from a number of linked components: a web server, a database management system and an application language. The architecture schema is shown in Figure 1.
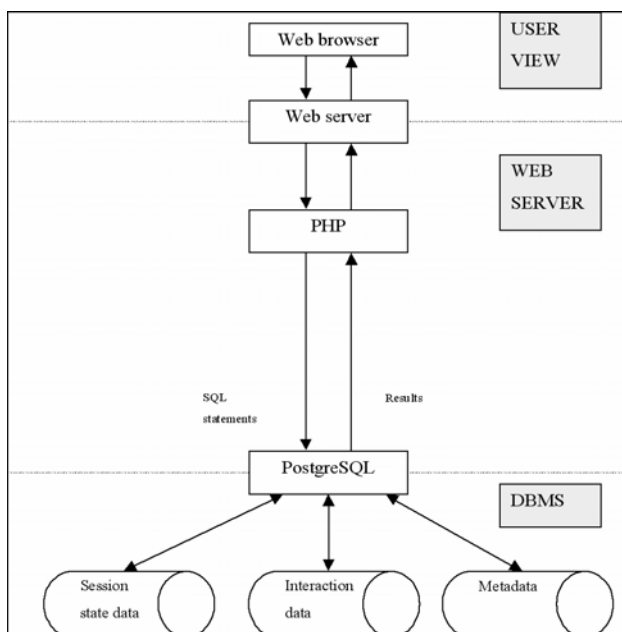


Figure 1: WICID architecture: schematic framework

Users interact with WICID by means of a web browser connected to a web server. WICID uses Apache (*http://www.apache.org*) as its web server. A database management system (DBMS) is used to store data of different types: primary migration and commuting flow data to which WICID provides access; metadata describing those primary

data; and 'state' data that record details of the sessions of each user. WICID uses PostgreSQL (*http://www.postgresql.org*) as its DBMS and to provide support for the storage and manipulation of geometric features (i.e. points, lines, polygons, etc.). In particular, a third party add-on to PostgreSQL called PostGIS (*http://postgis.refractions.net/*) offers facilities to handle spatial data that follow the OpenGIS "Simple Features Specification for SQL" standard. In order for dynamic web pages to be created, a programming language is required. WICID uses PHP (PHP Hypertext Processor) (*http://www.php.net*).

WICID has undergone considerable development over that last two years from the version that was reported in Stillwell and Duke-Williams (2004)[6]. In this section of the paper, we provide a short resume of the basic query-building procedure for users of the system before outlining two of the new features of the system, the map selection tool developed to support query-building and the analysis tool, designed to provide users with some insights into the data sets that they have extracted.

## Building queries in WICID

Since one of the fundamental aims of the CIDS has been to create a user-friendly interface, a great deal of effort has been directed at building a flexible, yet simple query interface so that queries can be formulated easily and data can be extracted and downloaded in a straightforward manner. It should be emphasised that users wanting to use the system will require an Athens userid and password and will have to have registered online to use the UK census data sets with the Census Registration Unit (http://census.data-archive.ac.uk/). Once the user has navigated to the homepage (http://cids.census.ac.uk/) and logged into the system using their Athens username and password, the user is confronted with the screen shown in Figure 2, containing a number of hotlinks that provide information about the data sets held in the system, details about the user's account and links to other useful web sites.
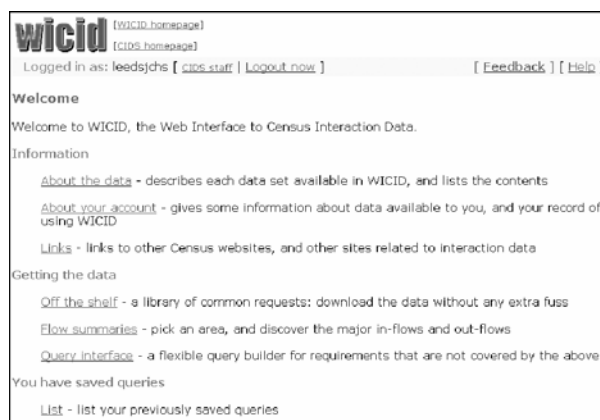


Figure 2: Welcome screen in WICID

There are three mechanisms for 'Getting the data' from WICID. The first is through the 'Off-the-shelf' facility of downloading data from a library of prepared queries. The second allows users to generate 'Flow summaries' for individual areas that they can specify. Figure 3 is an example of a summary in which the user has asked for commuting flows to and from the City of London in 1991, including flows from within that borough. The data comes from the 1991 SWS Set C and has been aggregated to show the top ten districts of origin (London

boroughs in this case) and of destination. The user can produce similar lists for other zones by clicking on any one of the other areas in either of the two lists; a flow pyramid can be generated by clicking on the pyramid icon in the left hand column of each list. Figure 4 illustrates an age-specific pyramid for commuters between Havering and the City of London.
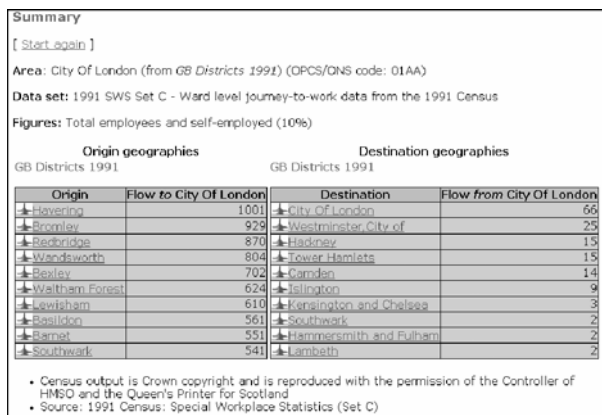


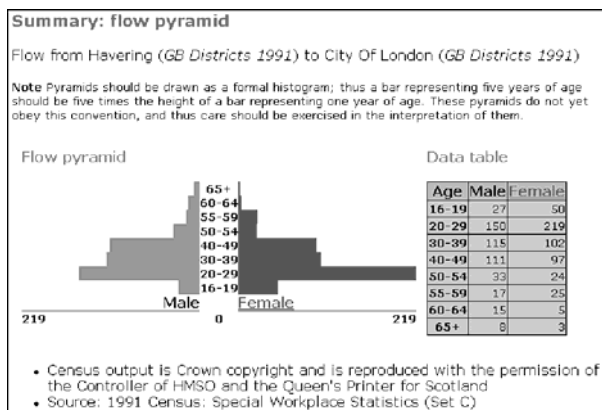Figure 3: Summary of commuting flows to and from the City of London, 1991



Figure 4: Age pyramid of commuters from Havering to the City of London, 1991

Building queries in WICID begins from the query interface (Figure 5) where the user is given the option of selecting 'Geography' or 'Data'. Of course, the interaction data sets all require double geographies and consequently the user has to decide which origin and destinations are required.
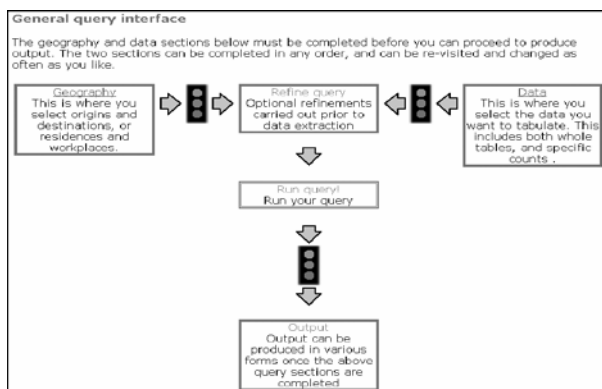


Figure 5: The general query interface in WICID prior to selection

One of the most innovative features of WICID is the facility to build geography selections in which the origins and destinations are drawn from sets of areas at different spatial scales and are not required to be the same. Consequently, it is possible, for example, for the user to select a single destination (e.g. the City of London) and to extract flows from other London boroughs, from other districts adjacent in the South East region and from other regions. If the flows required are commuters, then the user would make a selection of the appropriate variables from the table in the SWS.

There are various methods of selecting geographical areas: 'Quick selection' enables all areas at a certain scale to be selected; 'List selection' allows areas to be chosen from a list of all areas at each scale; and 'Type-in-box' selection provides for one area to be selected at a time. One of the new developments in WICID using web-mapping software is that of 'Map selection'. The map selection tool for choosing origins and destinations is crucial when users are unfamiliar with the geographical areas that they need to extract data for. We have found that this facility is particularly important when users are students doing project work, especially for small areas like wards.

The initial screen for map selection (Figure 6) contains three panels on the right hand side, the topmost of which is where the user chooses to select either areas or elements of areas by clicking the cursor in the map window. The panel below enables the user to zoom in and out or to reposition the map in the window. Finally, the lower panel contains various links to other screens where the user can change the map size, labelling and colour shading of the map and to reset the map to its original scale and view.
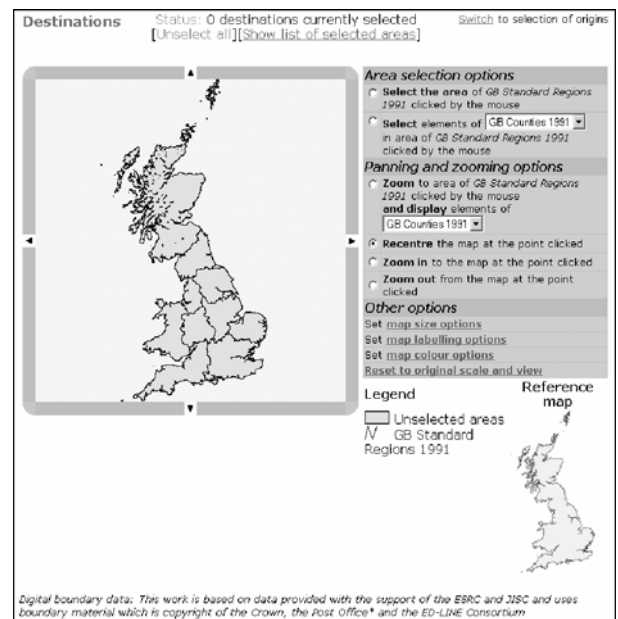


Figure 6: The map selection window in WICID

The example shown in Figure 7 illustrates the selection of the City of London at the district scale as a destination. Note that the status specification at the top of the figure changes as each destination is selected in the map window.
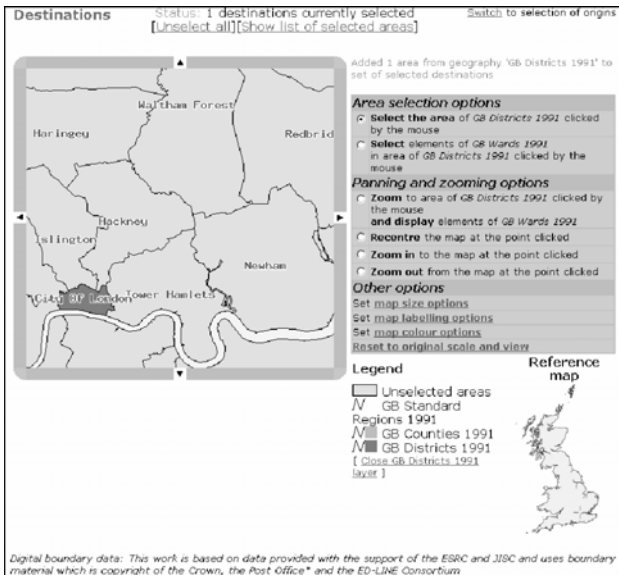
Figure 7: Example of map selection of City of London as destination area

Figure 8 is an example of a query to extract data on flows of employees and self-employed persons from the 2001 SWS for 110 origins (100 districts and 10 regions). Once the traffic lights are green, the user can run the query or refine the query in some way (e.g. by merging some of the variables selected) before extraction.



Figure 8: Query to select commuting data to City of London, 2001

**Output and analysis in WICID**

Once the selection has been made, the data can be extracted and downloaded. It is appropriate to check that everything is in order by returning to the general query screen and ensuring that the lights are on green (Figure 5) and then extracting the data by clicking on 'Run query'. The time taken to perform the query will depend upon the amount of data to be extracted. When the extraction has been completed the user will be informed of the time taken to complete the procedure (Figure 9) and invited to continue, at which point the data can be shown on screen (if not to extensive), downloaded or analysed. If the data is to be downloaded, the user has the option of choosing various formatting and layout options as indicated in Figure 10.
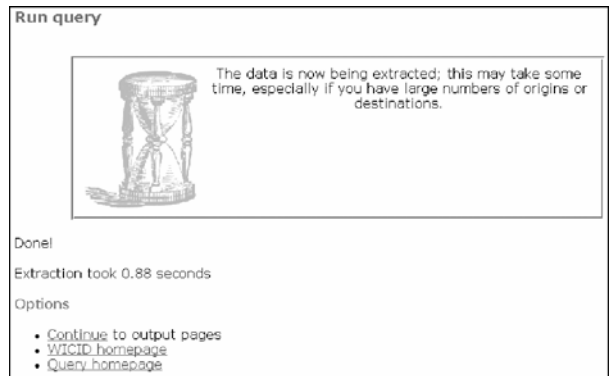


Figure 9: Screen indicating extraction has been completed
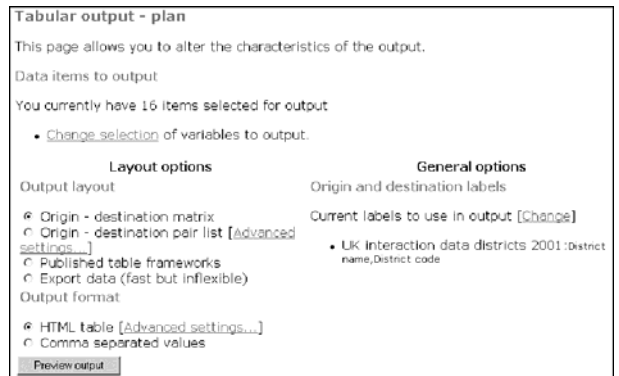


Figure 10: Output planning screen

However, WICID also allows users to undertake some analysis of the data extracted and 'add value' to the raw counts. The analytical facilities comprise a suite of five sets of indicators (Figure 11) that WICID will compute for a selection of any five of the counts that have been selected in the query. This maximum of five is to allow easy output on a single screen.
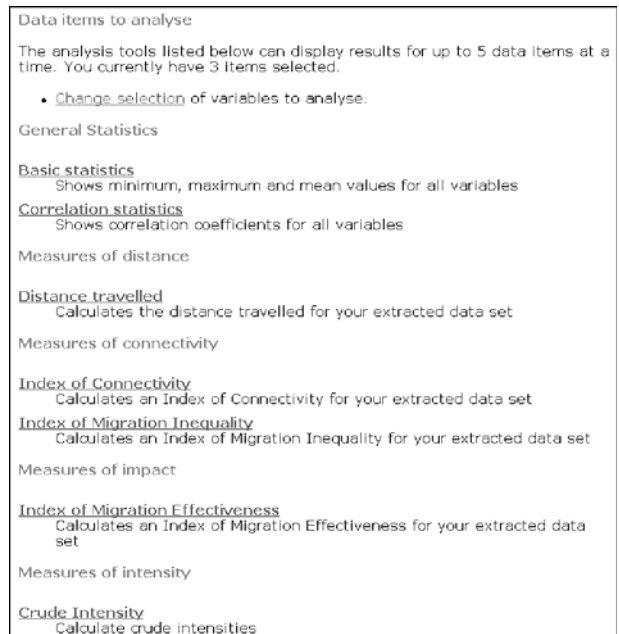


Figure 11: Analysis indicators available in WICID

The general statistics consist of descriptive statistics for all the flows extracted together with correlation coefficients that

indicate the strength of any statistical correlation between each pair of variables. The distance indicator is simply the average distance travelled and this measure relies on the availability of a additional data on distances between origins and destinations. The distance measure is not available when the origins and destination sets are drawn from different scales.

Additional data are also required for computation of the crude intensities. In most cases, these intensities are commuting or migration rates whose computation requires that the flow is divided by the appropriate population at risk (PAR). For some variable counts, the PARs are straightforward (e.g. the PAR for migration outflows of those in age group 1-4 is the population of the area aged 1-4 on census date obtained from the area statistics) but for other variable counts, the PARs are much less straightforward to define (e.g. PARs for outflows of moving groups) and may not be available from area statistics or standard tables.

Three further indicators are available in the current WICID system, none of which require additional information. These are indicators of connectivity, inequality and effectiveness. The latter two tend to be used for migration analysis whilst the index of connectivity, measured as the number of pairs of zones that have a flow between them divided by the total number of pairs of zones selected, can be used with commuting as well as migration data.

## 3. Mapping and analysis

By way of example, in this paper I have chosen to examine how commuting flows to London boroughs have changed between 1991 and 2001, focusing in particular on the changes in flows within London and those arriving from the rest of England and Wales. Aggregate statistics from the SWS for 1991 and 2001 suggest an overall increase of 20% in the commuting flows into London boroughs from 2.1 to 2.5 million persons per day. However, the increase in flows from outside London is much smaller (8.6%) than the increase in inflows from other London boroughs (25%). In fact the largest increase in commuting (nearly 60%) during the 1990s was taking place within boroughs.
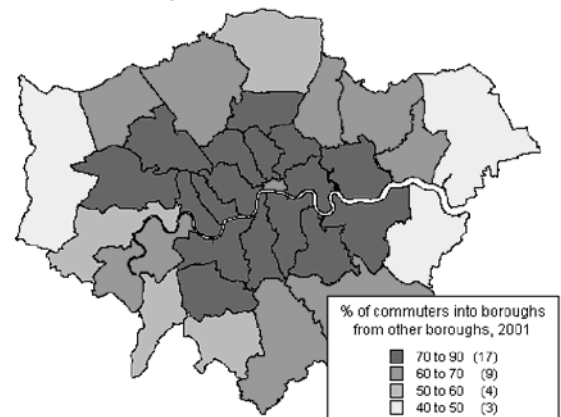
The variations across London between boroughs in terms of their commuting inflows are substantial but remain consistent between 1991 and 2001. The Borough of Westminster, for example, records over 350,000 commuters in 2001 with origins elsewhere in London and a further 110,000 coming from the rest of the country. The City of London ranks in second place with over 210,000 commuters from within London and further 100,000 from outside the capital. Only in two boroughs, Bexley and Havering, were the inflows from London exceeded by the inflows from elsewhere in both periods.

As we might expect, the geographical patterns of destinations for commuters depend to a large extent on where they originate. In Figure 12, two maps based on 2001 data are compared: map (a) illustrates that the inner London boroughs are the destinations for the large majority flows that originate within London, whilst the outer boroughs have much lower shares from within London. Map (b) is the mirror image of the former, demonstrating that for flows originating outside London, it is the outer boroughs (and the City of London) that have the
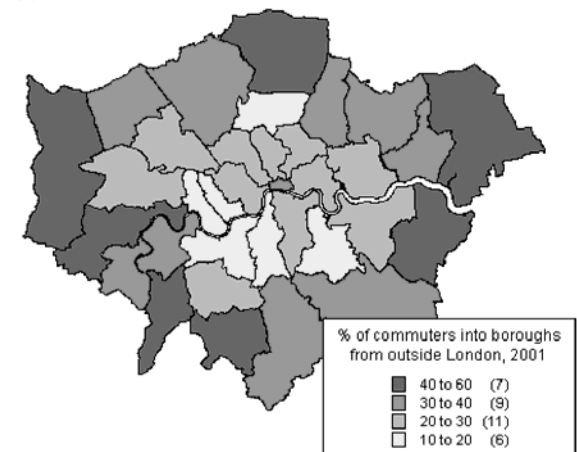
higher shares of this category of commuters. The patterns were much the same in 1991.



(a) London boroughs



(b) Inflows from elsewhere in London
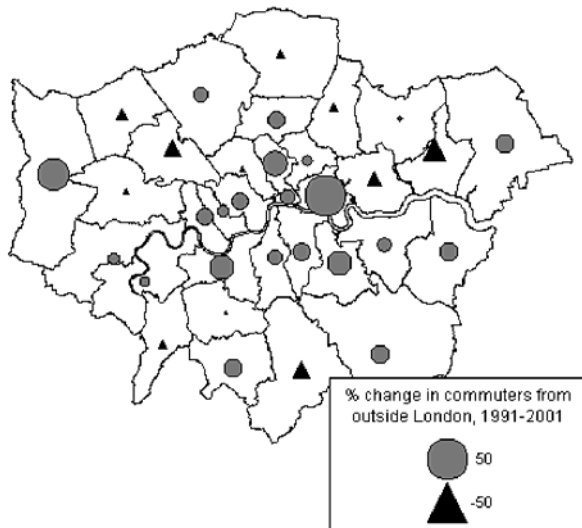


(c) Inflows from outside London

Figure 12: Percentage shares of commuting inflows to London boroughs, 2001

The increase in the overall volume of commuting between 1991 and 2001 conceals a range of experience for different boroughs (Figure 13). Commuting from within London to Tower Hamlets rose by over 70% whilst equivalent flows to Harrow declined by 14% between 1991 and 2001. Several boroughs experienced declines in their commuting inflows from outside London, particularly Brent, Croydon, Newham and Barking and

Dagenham, whilst the largest increases occurred in Tower Hamlets and Hillingdon.



(a) Inflows from elsewhere in London



(b) Inflows from outside London

Figure 13: Percentage changes in commuting inflows to London boroughs, 1991-2001

### 4. Conclusion

Since August 2006, CIDS has been relabelled as the Centre for Interaction Data Research and Estimation (CIDER), and will incorporate non-census flow data in future as well as interaction data from the 2011 Census. A user evaluation survey involving 165 respondents to an independent questionnaire indicated that almost 70% of users found the service to be either good or very good (Table 2). In 2004, users in 169 academic institutions carried out 3,246 sessions using WICID. However, use of the service increased significantly during the last quarter of 2004 and the first quarter of 2005 as the new 2001 interaction data sets have become available. We would expect the increase in usage to continue. Unfortunately the service is not available at the moment to academics or students outside the UK but it is hoped that this situation might be remedied in due course. An adapted version of the system (ACID) has been implemented that allows staff at the University of Queensland to have access to Australian internal migration data.

**Table 2: Results from 2004 user survey**

|  | Frequency | % Valid | % Cumulative |
|---|---|---|---|
| Very poor | 4 | 2.4 | 2.4 |
| Poor | 1 | .6 | 3.0 |
| Average | 46 | 27.9 | 30.9 |
| Good | 74 | 44.8 | 75.8 |
| Very Good | 40 | 24.2 | 100.0 |
| Total | 165 | 100.0 |  |

As well as the inclusion of new primary and derived data sets into the system, various functions or facilities within WICID have been developed to facilitate user access to the data prior to extraction and user analysis of the data following extraction. The map selection tool has been constructed in response to feedback from users (particularly students) whose knowledge of the geographies of the UK is limited and who require some assistance in selecting sets of zones to build customised queries involving origins and/or destinations at different spatial scales. The analysis facilities offered by WICID are useful because they add value to the raw counts extracted from the primary and derived data sets. Currently, a restricted set of indicators are available and considerable work remains to be done to assemble the count-specific populations at risk and the inter-area distances required to make the facilities fully functional with 2001 data. These tasks form part of the ongoing development work.

**References**
[1] P.J. Boyle and Z. Feng, "A method for integrating the 1981 and 1991 GB Census interaction data", **Computers, Environment and Urban Systems**, Vol. 26, 2002, pp. 241-56.
[2] K., Cole, M. Frost and F. Thomas, "Workplace data from the census", in P. Rees, D. Martin and P. Williamson, (eds.) **The Census Data System**, Chichester: Wiley, pp. 269-280, 2002.
[3] R. Flowerdew and A. Green, "Migration, transport and workplace statistics from the 1991 Census", in A. Dale and C. Marsh, (eds.) **The 1991 Census User's Guide**, London: HMSO, pp. 269-294, 1993.
[4] P.H. Rees and O.D. Duke-Williams, "Methods for estimating missing data on migrants in the 1991 British Census", **International Journal of Population Geography**, Vol. 3, , 1997, pp. 323-368.
[5] P.H. Rees, F. Thomas and O.D. Duke-Williams, "Migration data from the census", in P. Rees, D. Martin and P. Williamson, (eds.) **The Census Data System**, Chichester: Wiley, pp. 245-267, 2002.
[6] J.C.H. Stillwell, and O.D. Duke-Williams, "A new web-based interface to British census of population origin-destination statistics", **Environment and Planning A**, Vol. 35, 2003, pp. 113-132.